

# Practical - intermediate

Aurelien Ginolhac

2<sup>nd</sup> June 2016

## Project - set-up

- Create a new project in a meaningful folder name on your computer (such as `R_workshop/day1-intermediate`) using the project manager utility on the upper-right part of the rstudio window.
- Check if you have all those libraries installed

```
library("tidyr")
library("dplyr", warn.conflicts = FALSE)
library("ggplot2")
library("broom")
suppressPackageStartupMessages(library("GEOquery")) # bioconductor is verbose
theme_set(theme_bw(14)) # if you wish to get this theme by default
```

## Aim

Working with GEO datasets could be an hassle and you are going to experience it. Extensive manipulation of tables (`data.frame` and `matrix`) is required and provides a nice exercise. Here, we will investigate the relationship between the expression of *ENTPD5* and mir-182 as it was described by the authors. Even if the data are normalised and should be ready to use, quite an extensive amount of work is still required to reproduce the claimed result.

## Retrieve GEO study

The GEO dataset of interest is GSE35834

- load the study using the `getGEO` function

```
gse35834 <- getGEO("GSE35834", GSEMatrix = TRUE)
```

```
## ftp://ftp.ncbi.nlm.nih.gov/geo/series/GSE35nnn/GSE35834/matrix/
```

```
## Found 2 file(s)
```

```
## GSE35834-GPL15236_series_matrix.txt.gz
```

```
## Using locally cached version: /var/folders/7x/14cplkhj0fn34yltb3c0j9bczm4jkt/T//RtmpX6KADI/GSE35834-
```

```
## Using locally cached version of GPL15236 found here:
```

```
## /var/folders/7x/14cplkhj0fn34yltb3c0j9bczm4jkt/T//RtmpX6KADI/GPL15236.soft
```

```
## Warning in read.table(file = file, header = header, sep = sep, quote =
```

```
## quote, : not all columns named in 'colClasses' exist
```

```
## GSE35834-GPL8786_series_matrix.txt.gz
```

```
## Using locally cached version: /var/folders/7x/14cplkhj0fn34yltb3c0j9bczm4jkt/T//RtmpX6KADI/GSE35834-
```

```
## Using locally cached version of GPL8786 found here:
```

```
## /var/folders/7x/14cplkhj0fn34yltb3c0j9bczm4jkt/T//RtmpX6KADI/GPL8786.soft
```

```
## Warning in read.table(file = file, header = header, sep = sep, quote =  
## quote, : not all columns named in 'colClasses' exist
```

```
show(gse35834)
```

```
## `$`GSE35834-GPL15236_series_matrix.txt.gz`  
## ExpressionSet (storageMode: lockedEnvironment)  
## assayData: 22486 features, 80 samples  
## element names: exprs  
## protocolData: none  
## phenoData  
## sampleNames: GSM875933 GSM875934 ... GSM876012 (80 total)  
## varLabels: title geo_accession ... data_row_count (39 total)  
## varMetadata: labelDescription  
## featureData  
## featureNames: 10000_at 10001_at ... 9_at (22486 total)  
## fvarLabels: ID ENTREZ_GENE_ID Description SPOT_ID  
## fvarMetadata: Column Description labelDescription  
## experimentData: use 'experimentData(object)'  
## Annotation: GPL15236  
##  
## `$`GSE35834-GPL8786_series_matrix.txt.gz`  
## ExpressionSet (storageMode: lockedEnvironment)  
## assayData: 7815 features, 78 samples  
## element names: exprs  
## protocolData: none  
## phenoData  
## sampleNames: GSM875855 GSM875856 ... GSM875932 (78 total)  
## varLabels: title geo_accession ... data_row_count (39 total)  
## varMetadata: labelDescription  
## featureData  
## featureNames: 14q-0_st 14qI-1_st ... zma-miR408_st (7815 total)  
## fvarLabels: ID miRNA_ID_LIST ... SEQUENCE (11 total)  
## fvarMetadata: Column Description labelDescription  
## experimentData: use 'experimentData(object)'  
## Annotation: GPL8786
```

- what kind of object is `gse35834`?
- Two platforms were used in this study, which ones?
- How can you assign the mRNA or mir data to each element of `gse35834`?

## Explore the mRNA expression meta-data

Informations about samples are accessible using `phenoData()` and can directly be retrieved as a `data.frame` with `pData()`.

for example, the following command will return the mRNA meta-data as a `data.frame`

```
pData(gse35834[[1]])
```

- Extract as a `tbl_df` named `rna_meta` the mRNA meta-data and
  - rename `geo_accession` to `sample`
  - select `source_name_ch1` and all columns that start with “charact”

## Explore the mir expression meta-data

- Extract as a `tbl_df` named `mir_meta` the mir meta-data and
  - rename `geo_accession` to `sample`
  - select `source_name_ch1` and all columns that start with “charact”

## Join meta-data

- Explore the two data frames with `View(rna_meta)` and `View(mir_meta)`. Are the samples `GSM*` identical?

Then, we would like to somehow join both informations.

Knowing that both data frames have different “sample” columns, merge them to get the correspondence between RNA `GSM*` and mir `GSM*`. Save the result as `rna_mir`.

### Note

When 2 data.frames are joined by specific columns and the remaining columns have identical names, a ‘x’ or ‘y’ suffix is appended for the first and second data frames respectively

## Get RNA expression data for the *ENTPD5* gene

Expression data can be accessed via `exprs()` which returns a matrix.

### Warning

If you do not pipe the command to `head`, R would print **ALL** rows (or until it reaches `max.print`).

```
exprs(gse35834[[1]]) %>% head()
```

rows are probes and columns are sample ids in the form `GSM*`.

Probe ids are not meaningful, but `fData()` provides features.

```
fData(gse35834[[1]]) %>% head()
```

Again, we need to merge both informations to assign the expression data to the gene of interest.

1. Find the common values that could help us joining.
2. A `matrix` contains only numerical values. But, the `rownames` contain the necessary info. Transform the `matrix` into a `data.frame`. Then, convert the `rownames` to a column using `tibble::rownames_to_column(var = "ID")`. Save as `rna_expression`
3. merge expression data to platform annotation (`fData(gse35834[[1]])`). Save as `rna_expression`. R is always working on temporary objects, you won’t erase the object you are working on.

### Note

Warnings about *factors being coerced to characters* can be ignored. Factors shouldn’t be in the first place (default of `readr` functions)

4. Find the Entrez gene id for *ENTPD5*. Usually, the gene symbol is given in the annotation, but each GEO submission is a new discovery.
5. Filter `rna_expression` for the gene of interest and tidy the samples:  
A column `sample` for all `GSM*` and a column `rna_expression` containing the expression values. Save the result as `rna_expression_melt`. At this point you should get a `data.frame` of 80 values.

6. Add the meta-data and discard the columns `ID`, `SPOT_ID` and `sample.x`. Save the result as `rna_expression_melt`.

## Get mir expression data for miR-182

1. Repeat the previous step but using `exprs(gse35834[[2]])` for the `mir_expression`. This time, the mir names are nicely provided by `fData(gse35834[[2]])` in the column `miRNA_ID_LIST`
2. How many rows do you obtain? How many are expected?
3. Find out what happened, and plot the boxplot distribution of `expression` by `ID`
4. Filter out the irrelevant IDs using `grepl` in the `filter` function.

### Hint

adding `!` to a condition means NOT. Example `filter(iris, !grepl("a", Species))`: remove all `Species` that contain an "a".

5. Add the meta-data, count the number of rows. Discard the column `sample.x` after joining.

## join both expression

Join `rna_expression_melt` and `mir_expression_melt` by their common columns EXCEPT `sample`. Save the result as `expression`.

## Examine gene expression according to meta data

1. Plot the gene expression distribution by Gender. Is there any obvious difference?
2. Plot gene AND mir expression distribution by Gender. Is there any obvious difference?

### Hint

You will need to tidy by gathering rna and mir expression

3. Plot gene AND mir expression distributions by source (control / cancer). To make it easier, a quick hack is `separate(expression, source_name_ch1, c("source", "rest"), sep = 12)` to get `source` as control / cancer. Is there any difference?
4. Replot 3. but reordering the levels so normal colon comes first. Display *normal* in "lightgreen" and *cancer* in "red" using `scale_fill_manual()`

## plot relation ENTPD5 ~ mir-182 as scatter-plot for all patients

- add a linear trend using `geom_smooth()` for all data + per source
- does it support the claim of the study?

## Supplementary exercise - linear regression

- get the estimate from the linear trend. linear models are outputted by `lm()` as lists. Since `data.frame` are much easier to work with, David Robinson developed `broom`. We will present the use of `broom` during the advanced lecture, but you can have an insight here and test `broom::tidy()` coupled with `dplyr::do()`.

```

library("broom")
expression %>%
  separate(source_name_ch1, c("source", "rest"), sep = 12) %>%
  group_by(source) %>%
  do(tidy(lm(rna_expression ~ mir_expression, data = .))) %>%
  filter(term != "(Intercept)")

```

```

## Source: local data frame [2 x 6]
## Groups: source [2]
##
##           source      term  estimate  std.error  statistic  p.value
##          (chr)      (chr)    (dbl)     (dbl)      (dbl)     (dbl)
## 1 colon cancer mir_expression -0.03354124 0.09217281 -0.3638951 0.7174117
## 2 normal colon mir_expression  0.04496954 0.10042656  0.4477853 0.6588934

```

The estimate of the intercept is not meaningful thus it is filtered out. One can easily see that the slope is not significant when data are slipped by source.

- Perform the linear regression and tidy the results for all data, is it significant?
- replace `tidy` by `glance` to extract the  $r^2$ . Is this value satisfactory?

## Perform a linear model for the expression of *ENTPD5* and ALL mirs

- Count how many `hsa-mir`, which are not star, are present on the array GPL8786
- Retrieve the expression values for the 677 human mir like you did before. Same procedure, except that you don't filter for mir-182. Save as `all_mir_rna_expression`
- Perform the 677 linear models, tidy the results and arrange by the `adj.r.squared`
- Get the top 12 mir and plot the scatter plot